# A Model for Saliency-Based Visual Attention for Rapid Scene Analysis

Laurent Itti, Christof Koch and Ernst Niebur

Jyothsna Shashikumar Sastry

# Overview

- Motivation
  - Introduction to Saliency based Model
- Architecture
  - Extraction of Early Visual Features
  - Saliency Maps
- Comparison with Spatial Frequency Content Models
- Results
- Strengths and limitations
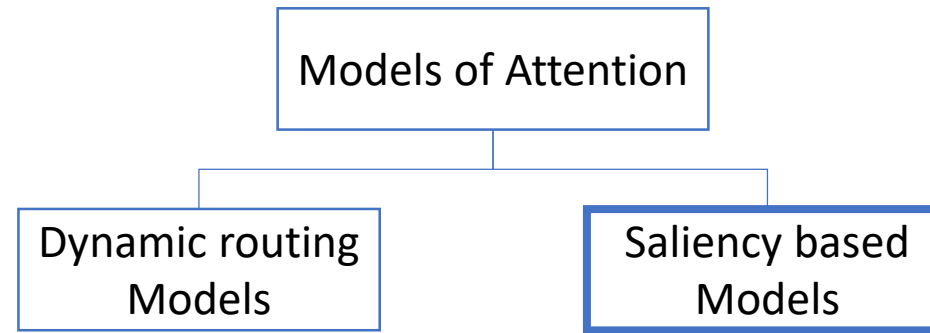- Summary

# Activity

## Analyse this scene

# MOTIVATION

- What did you see? List the salient objects you identified.

- We can interpret complex scenes in real time

- Focus of Attention (FOA) processes only a subset of the sensory information to reduce complexity of scene analysis

**GOAL : Build a model that mimics primate visual attention for static scenes**
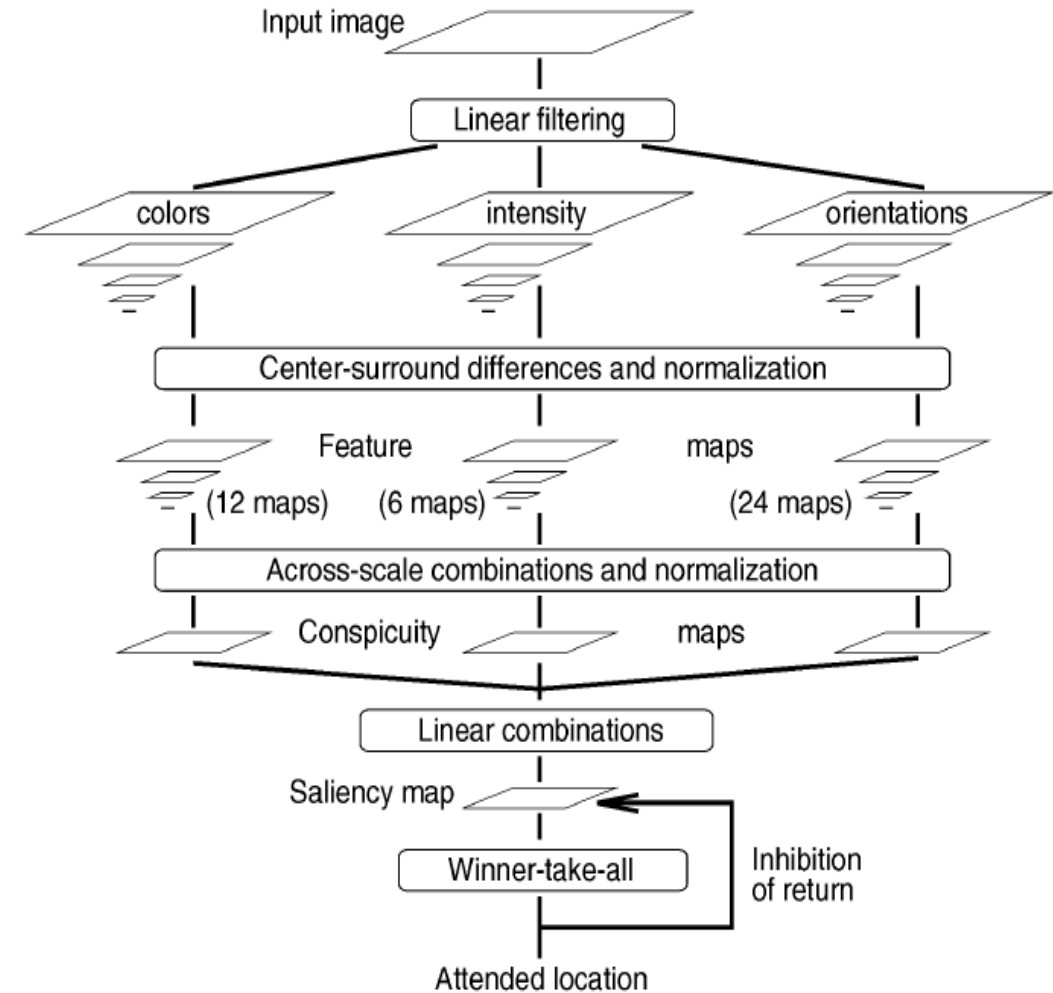
```
                    ┌─────────────────────┐
                    │ Models of Attention │
                    └─────────┬───────────┘
                   ┌──────────┴──────────┐
          ┌────────────────┐    ┌────────────────┐
          │ Dynamic routing│    │ Saliency based │
          │     Models     │    │     Models     │
          └────────────────┘    └────────────────┘
```

- Based on feature integration theory

- Visual input is decomposed into a set of feature maps

- Different spatial locations within a map compete for saliency

- All feature maps feed into a master "saliency map"

Massively parallel method for the rapid selection of a small number of interesting image locations to be analysed by more complex and time-consuming object-recognition processes

# ARCHITECTURE

- Input:  Static color images digitized at 640 X 480
- Multiscale feature extraction: Nine spatial scales using Gaussian pyramids ( σ ∈ [0..8] where σ is the scale )
- **Center surround** implemented to detect locations that stand out from their surround
  - **Center**: a pixel at scale c ∈ {2, 3, 4}
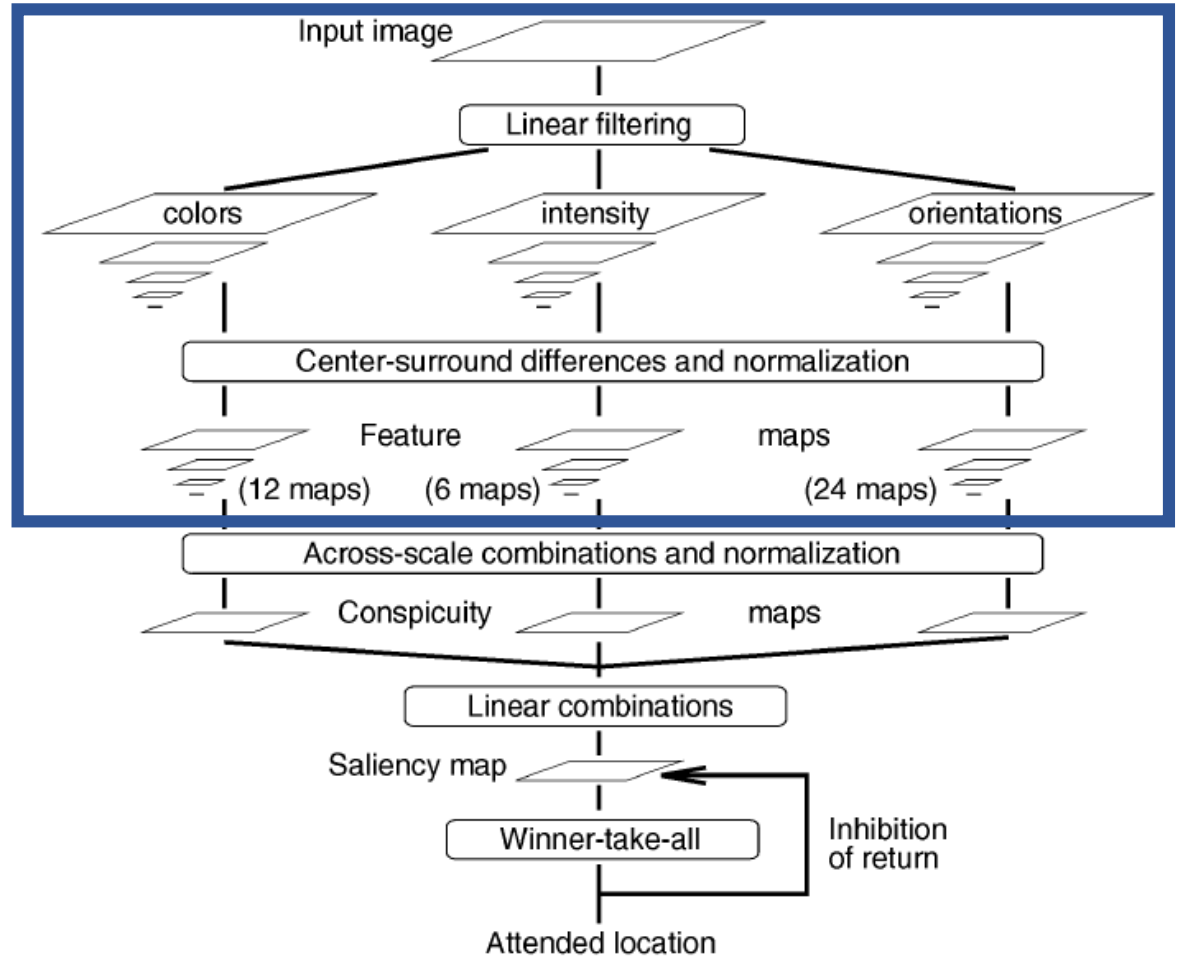  - **Surround** is the corresponding pixel at scale s = c + δ, with  δ ∈ {3, 4}.

The across-scale difference between two maps, denoted "⊖"

# Extraction of Early Visual Features

**Linear Filtering** and **Center-surround differences and Normalization**

- Intensity
- Color
- Orientation

# Extraction of Early Visual Features

**Intensity contrast maps**

- Mammals are equipped with neurons sensitive to light centers with dark surrounds and dark centers with light surrounds
- Intensity image $I = (r + g + b) / 3$
- Gaussian pyramid I(σ) where σ ∈ [0..8]
- Normalize each channel by I only where I > (1/10)th of $I_{max}$ (maximum I over entire image)
- Intensity contrast maps $I(c, s) = |I(c) \ominus I(s)|$

  c ∈ {2, 3, 4} and s = c + δ, with  δ ∈ {3, 4}.

# Extraction of Early Visual Features

**Color feature maps**

- "Color double-opponent" system in the human cortex
- Four broadly-tuned color channels are created
  - $R = r - (g + b)/2$,
  - $G = g - (r + b)/2$
  - $B = b - (r + g)/2$
  - $Y = (r + g)/2 - |r - g|/2 - b$
- Four Gaussian Pyramids from these channels
  - $R(\sigma)$, $G(\sigma)$, $B(\sigma)$ and $Y(\sigma)$
- According to the color opponency,

$RG(c, s) = |(R(c) - G(c)) \ominus (G(s) - R(s))|$

$BY(c, s) = |(B(c) - Y(c)) \ominus (Y(s) - B(s))|$

# Extraction of Early Visual Features

**Orientation maps**

- Human visual attention system has the ability to identify breaks in patterns with the aid of orientation sensitive neurons of the cortex
- Intensity image **I** convolved with an Orientational filter (Gabor filters)
- $\Theta \in \{0^0, 45^0, 90^0, 135^0\}$
- Gabor pyramids at nine scales [0...8]

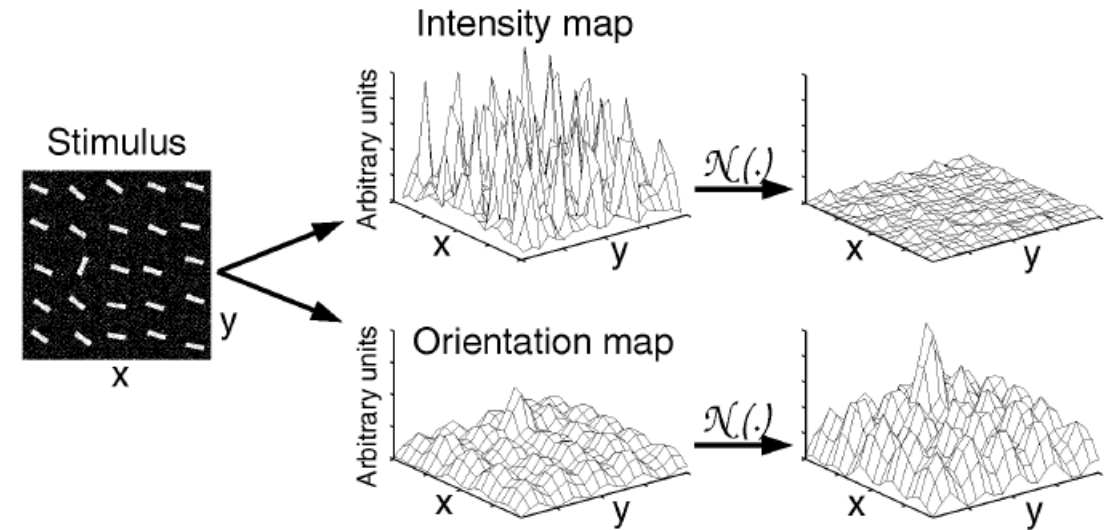$\boldsymbol{O}(c, s, \Theta) = |O(c, \Theta) \ominus O(s, \Theta)|$

# Saliency Maps

- 42 Feature maps:
  - 6 for Intensity
  - 12 for color
  - 24 for orientation
- Difficulty in combining different feature maps?

# Saliency Maps

- Normalization operator **N**(.)
  - Normalize map to a fixed range [0…M]
  - Find location of global maximum M
  - Compute average $\bar{m}$ of the local maxima
  - Multiply map by $(M-\bar{m})^2$

# Saliency Maps

- Feature maps are combined into three "conspicuity maps" by across-scale additions at σ = 4 of the saliency map

$$\overline{I} = \overset{4}{\underset{c=2}{\oplus}} \overset{c=4}{\underset{s=c+3}{\oplus}} \mathcal{N}\left(I(c,s)\right)$$

$$\overline{C} = \overset{4}{\underset{c=2}{\oplus}} \overset{c+4}{\underset{s=c+3}{\oplus}} \left[\mathcal{N}\left(\mathcal{R}\mathcal{G}(c,s)\right) + \mathcal{N}\left(\mathcal{B}\mathcal{Y}(c,s)\right)\right]$$

$$\overline{O} = \underset{\theta\in\{0°,45°,90°,135°\}}{\sum} \mathcal{N}\left(\overset{4}{\underset{c=2}{\oplus}} \overset{c+4}{\underset{s=c+3}{\oplus}} \mathcal{N}\left(O(c,s,\theta)\right)\right)$$

- The three conspicuity maps are averaged into the final input **S** to the saliency map

$$\mathbf{S} = \left(\mathbf{N}(\overline{I}) + \mathbf{N}(\overline{C}) + \mathbf{N}(\overline{O})\right) / 3$$

# Saliency Maps

Model the Saliency Map as a 2D layer of *integrate and fire* neurons at scale 4
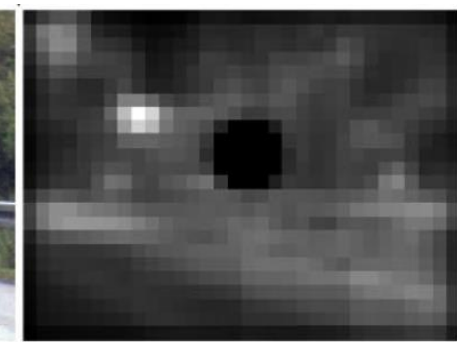
Working:

- Neurons in the SM receive excitatory inputs from *S*
- Each SM neuron excites its corresponding Winner Take All neuron.
- The first "winner" to reach a threshold fires.
- FOA shifted to this location
- All WTA neurons are reset
- SM neurons are reset in the location of FOA – aiding next salient location to become the winner
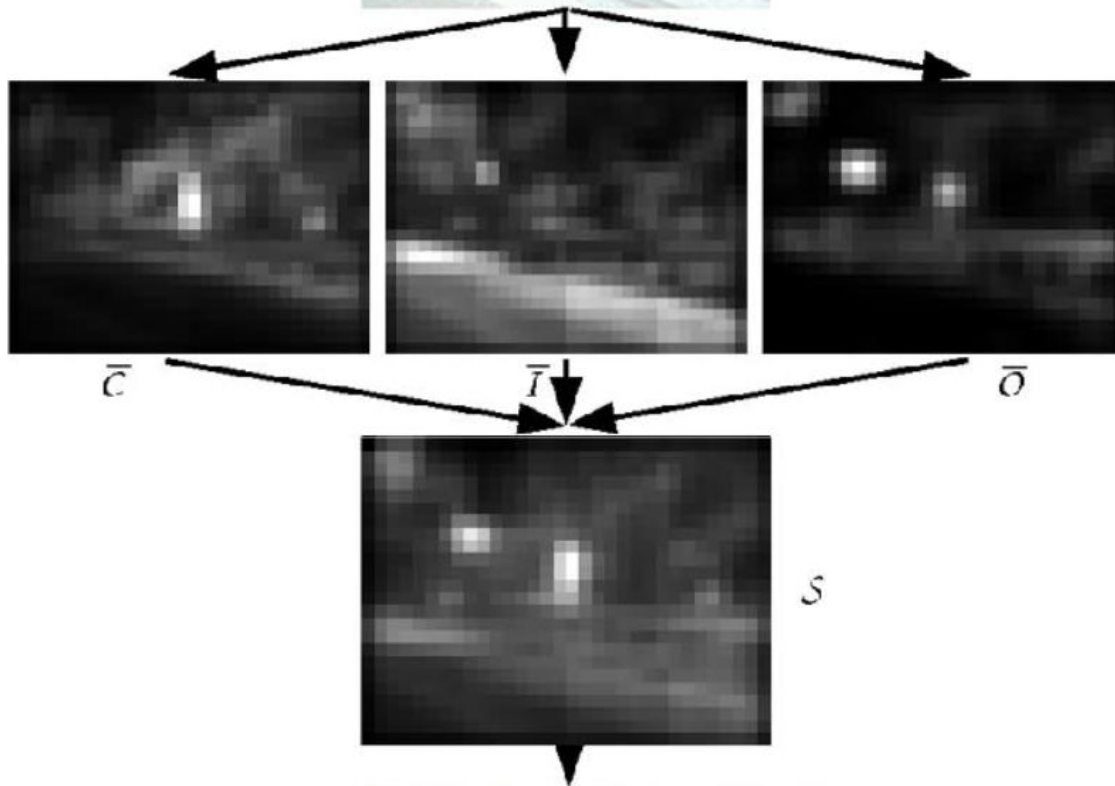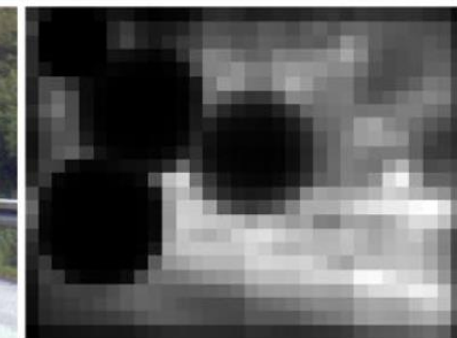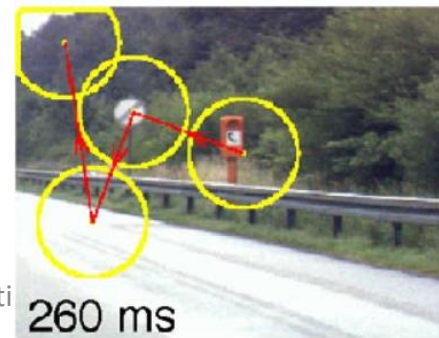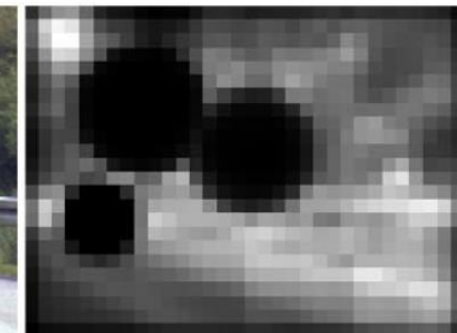
Input image

Focus of Attention

Saliency maps

$\bar{C}$    $\bar{I}$    $\bar{O}$

$S$

92 ms

145 ms

206 ms

260 ms

# Comparison with Spatial Frequency Content Models
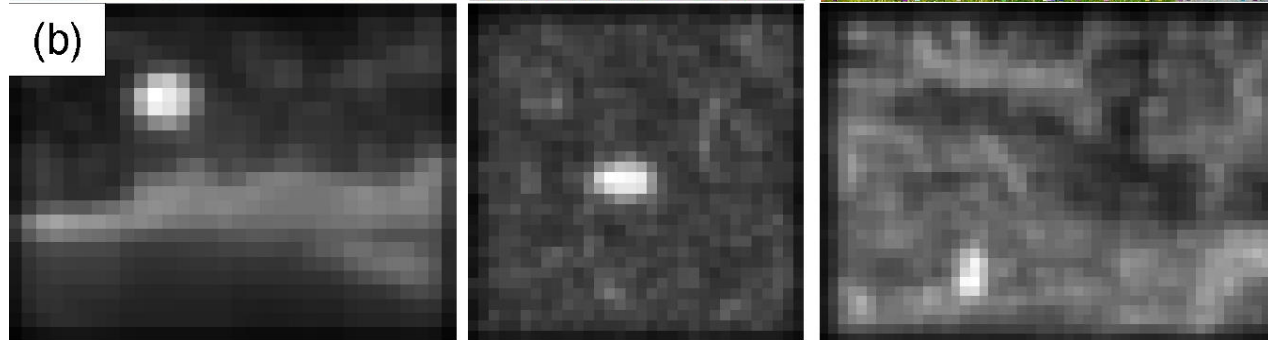
- At a given image location, a 16 X 16 image patch is extracted from each I(2), R(2), G(2), B(2), and Y(2) map
- 2D Fast Fourier Transforms (FFTs) are applied to the patches
- The SFC measure is the average of the numbers of nonnegligible coefficients in the five corresponding patches

- Results show that our model is superior to spatial frequency content (SFC) based models in the presence of noise
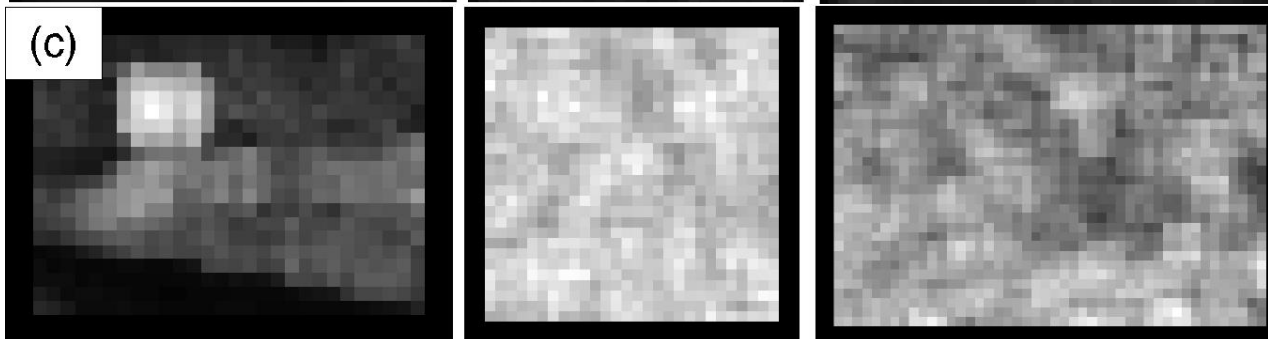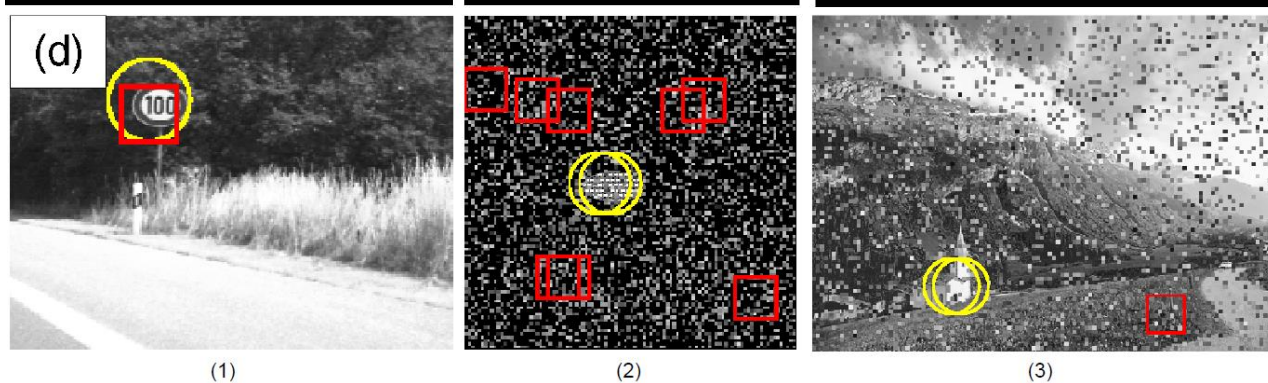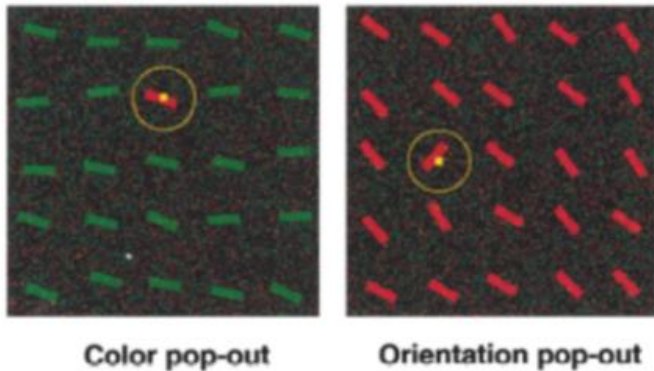
Input Images

Saliency maps

SFC maps

Output salient
Locations
SM(Yellow circle)
SFC(Red square)

(1)　　　　　　(2)　　　　　　(3)

17

# RESULTS

## General Performance

- Extensively tested with artificial images to ensure proper functioning
- Robust to addition of noise
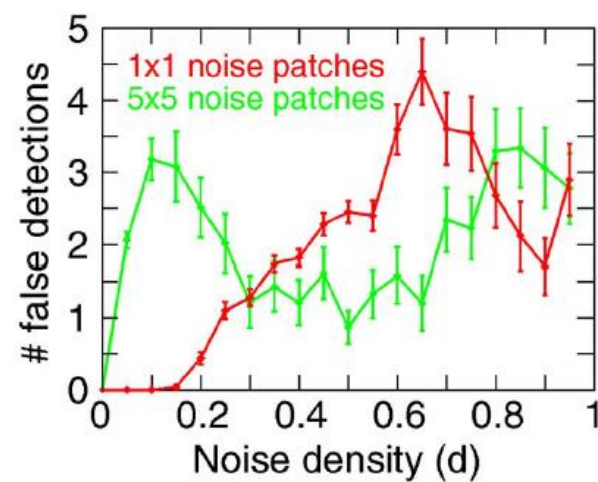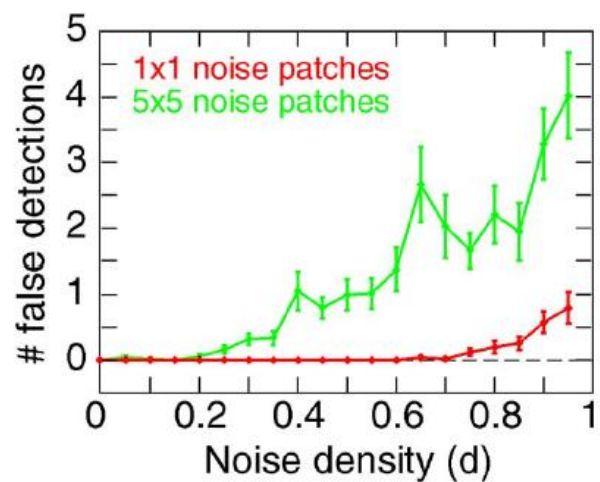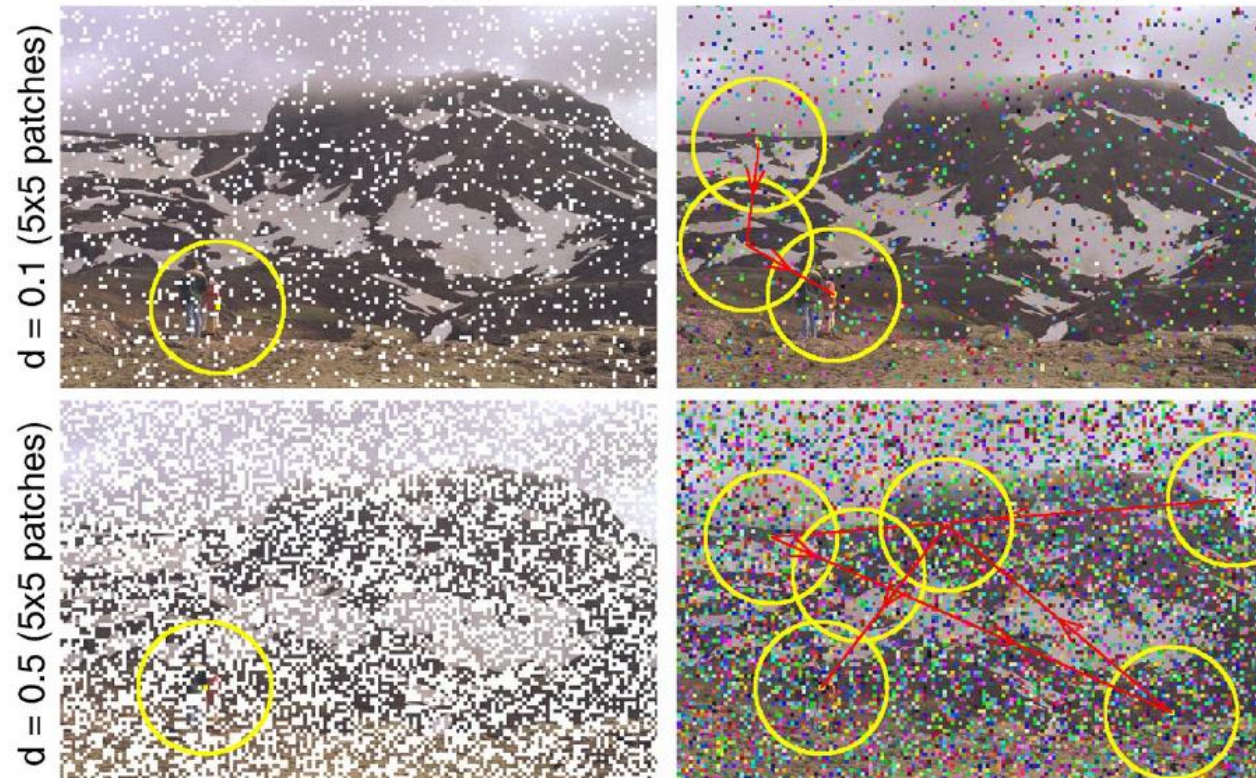- Reproduces human performance for pop-out tasks



Color pop-out    Orientation pop-out

- Tested with real images: attended locations were the objects of interest

A Model of Saliency-Based Visual Attention for Rapid Scene Analysis

# Strengths

✓ Mimics the properties of primate vision

✓ Despite its simple architecture and feed-forward feature-extraction mechanisms, the model is capable of strong performance with complex natural scenes

✓ Massively parallel implementation for the rapid selection of a small number of interesting image locations

✓ Allows real time operation

# Limitations

- Only object features explicitly represented in at least one of the feature maps can lead to pop-out

- Fails to detect targets salient for unimplemented feature types (e.g., T junctions or line terminators)

- Cannot reproduce contour completion or closure

- No magnocellular motion channel

# SUMMARY

- A conceptually simple computational model for saliency-driven visual attention.

- Architecture is inspired by biological insights

- Efficient in reproducing some of the performances of primate visual systems.

- The efficiency of this approach for target detection critically depends on the feature types implemented.

- Model can be easily tailored to arbitrary tasks using dedicated feature maps.

A Model of Saliency-Based Visual Attention for Rapid Scene Analysis

# REFERENCES

- **A Model of Saliency-Based Visual Attention for Rapid Scene Analysis**
  Laurent Itti, Christof Koch, and Ernst Niebur

- https://github.com/mbanani/attend

- Image source: Google

# Thank you!

A Model of Saliency-Based Visual Attention for Rapid Scene Analysis